

# PATRIC update

BRC3

February 2006

# PATRIC's Pathogens

- Bacteria:
  - *Brucella*
  - *Rickettsia*
  - *Coxiella burnetii*
- Viruses:
  - Coronaviruses (SARS)
  - Caliciviruses
  - Hepatitis A viruses
  - Hepatitis E viruses
  - **Rabies viruses** ⇒ **Lyssaviruses** (expanded scope)

# Curation Concepts

- **Nucleotide-level Curation**
  - CDSs and RNA gene calls
  - Ribosome binding sites and start site correction
  - Pseudogenes and other physical features
- **Protein-level Curation**
  - Functional assignment, classification
  - Structural properties and features
- **Automated and Manual Curation**
- **Reference and Associated Genomes**

# Curation Progress

- All bacterial and viral **reference genomes** have received **nucleotide-level manual** curation
  - latest release: **Dec 22, 2005**
- **Protein-level automated** curation has been done on all bacterial reference genomes (not public yet)

# Summary of Nucleotide-level Curation for Reference Genomes

Pathosystem	#genes GenBank entry	Start-site changes	Frameshifts	Premature stops	Genes deleted	Genes added
<i>Brucella</i>	3198	961	100	50	6	350
<i>Coxiella</i>	2052	74	1	0	0	124
<i>Rickettsia</i>	835	32	22	1	0	74

	No. of RGs	Total genes/genome	1st pass check	New annotations
Calicivirus	13	2,3	Y	
Coronavirus	16	14	Y	SARS (SZ-3)
Hepatitis A	1	2	Y	
Hepatitis E	5	3,4	Y	
Lyssavirus	1	5	Y	

# Literature Curation: PathInfo

**Brucella melitensis**

**I. Organism Information**

**A. Taxonomy Information**

1. **Species:**

- a. *Brucella melitensis* (Morenoa et al., 2002):
  - i. GenBank Taxonomy No.: 29459
  - ii. Description: According to the new taxonomy used by NCBI, the classic *Brucella* spp. are named *Brucella melitensis*, which include *Brucella melitensis* 16M and 5 biovars: Abortus, Canis, Neotomae, Ovis, and Suis. However, the traditional taxonomy is still widely used, where *Brucella* spp. include *B. melitensis*, *B. abortus*, *B. suis*, *B. canis*, *B. neotomae*, and *B. ovis*. This document is focused on the first three species.
  - iii. Variant(s):
    - *Brucella melitensis* biovar 1 strain 16M (Morenoa et al., 2002, Gandara et al., 2001):
      - o GenBank Taxonomy No.: 224914
      - o Parent: *Brucella melitensis*
      - o Description: Strain 16M, corresponding to ATCC 23456, is the type strain for this biovar. Strain 16M primarily affects goats and sheep, and is the most virulent of the *Brucellae* in humans.
    - *Brucella melitensis* biovar Abortus (Morenoa et al., 2002):
      - o GenBank Taxonomy No.: 235
      - o Parent: *Brucella melitensis*

References:

Morenoa, 1999; Morenoa, F. The Epidemiology of Bovine Brucellosis. *Advances in Veterinary Science and Comparative Medicine*. 1980; 24: 69 - 98. [PubMed: 6779513].

Orduna et al., 2000: Orduna A, Almaraz A, Prado A, Gutierrez MP, Garcia-Pascual A, Duenas A, Cuervo M, Abad R, Hernandez B, Lorenzo B, Bratos MA, Rodriguez-Torres A Evaluation of an immunature-acclimation test (Brucellacart) for serodiagnosis of human brucellosis

## Curated Information

- Organism background
- Epidemiology
- Pathogenesis
- Experiments

Pathosystem	
Bacteria	Status
<i>Brucella</i>	Available
<i>Coxiella</i>	Available
<i>Rickettsia</i>	Available
Viruses	
Calicivirus	Ready, awaiting release
Coronavirus	in queue
Hepatitis A	in queue
Hepatitis E	Ready, awaiting release
Lyssavirus	In progress

# PathInfo Status for the Pathogens of Other BRCs

<b>Eukaryotic pathogens</b>		
<b>BRC</b>	<b>complete</b>	<b>to be completed</b>
Apicomplexan	2	1
BioHealthBase	0	1
Pathema	0	1
<b>Viruses</b>		
	<b>complete</b>	<b>to be completed</b>
VBRC	16	4
BioHealthBase	0	1
<b>Bacteria</b>		
	<b>complete</b>	<b>to be completed</b>
Pathema	3	2
NMPDR	0	4
ERIC	4	1
BioHealthBase	1	0



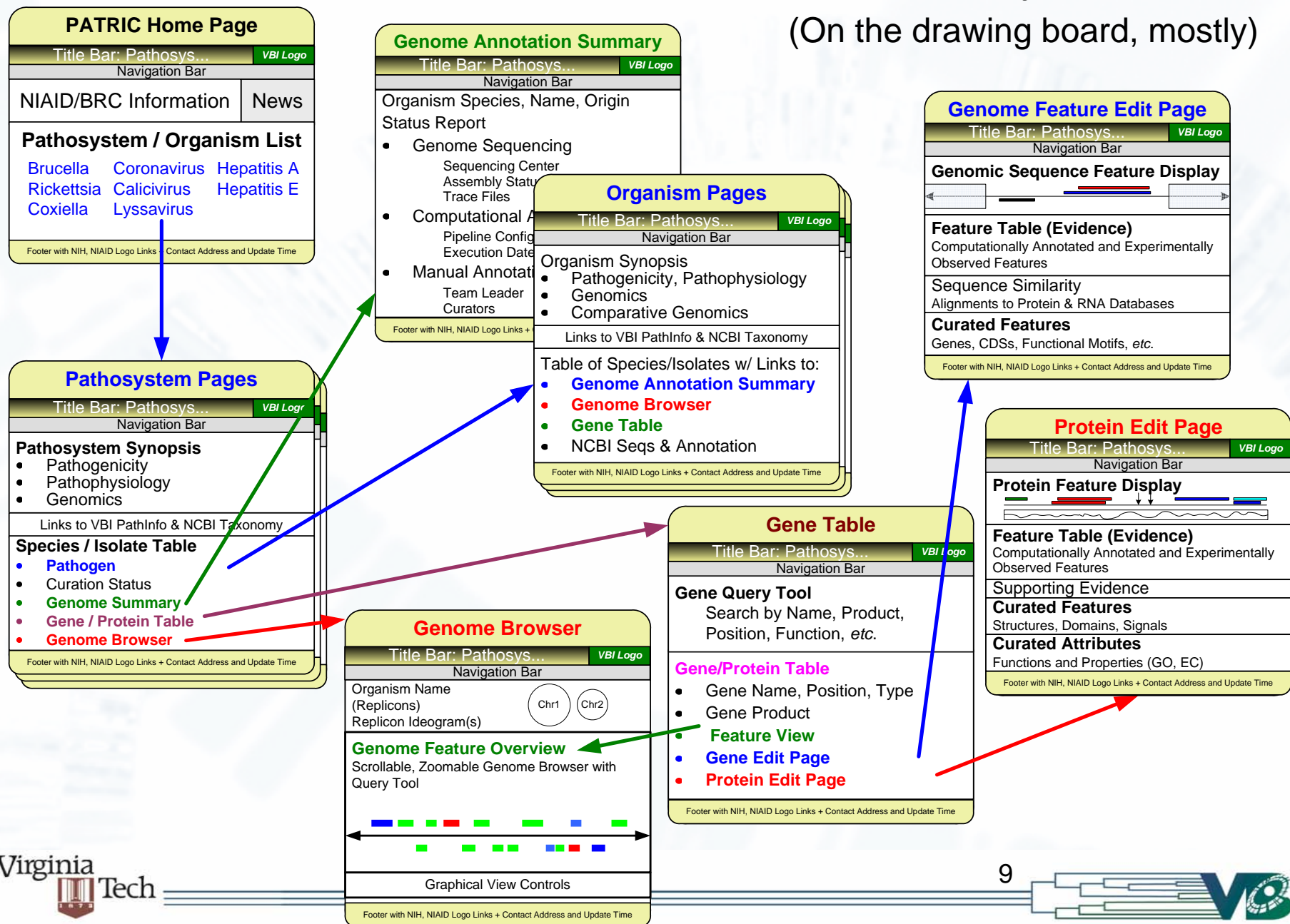
# Progress: Software Development

- **Curation Infrastructure**
  - Nucleotide-level curation: pipeline + edit page
  - Protein-level curation: pipeline + edit page
- **Improvements to External Website**
  - User-friendly database query interface
  - Enhanced gene/protein table display highlighting corrections to published annotations

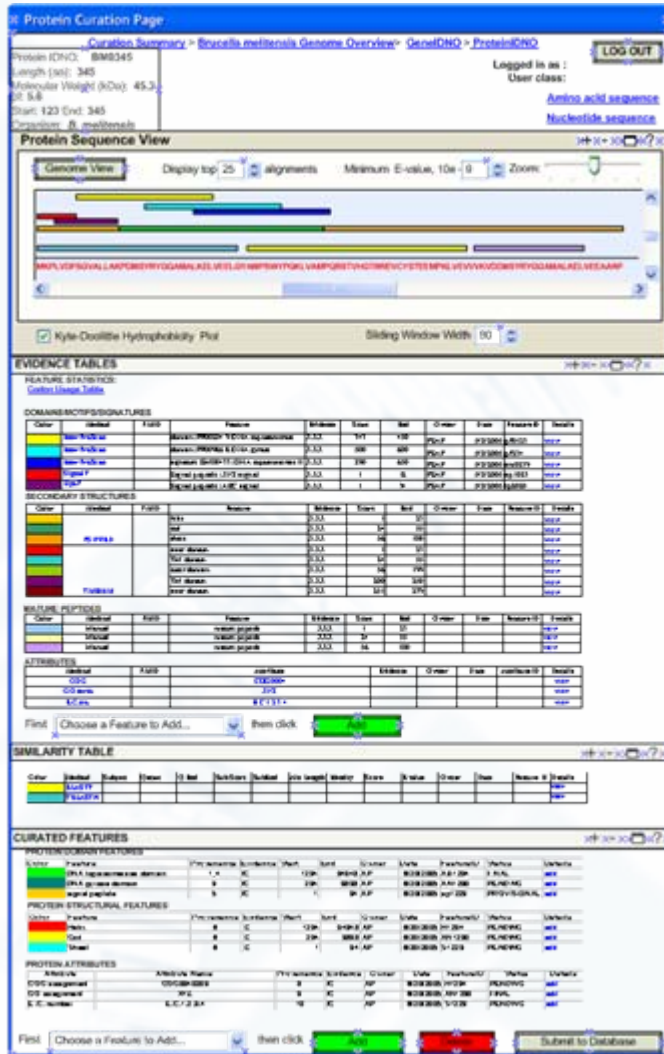


# Curation Infrastructure: February '05

(On the drawing board, mostly)

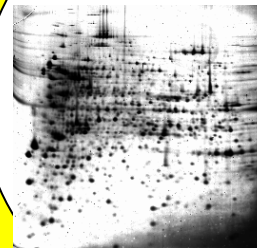




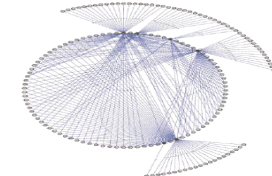


## In the works: links from protein edit page to proteomics data

2D-gel

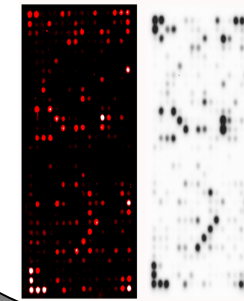


Y2H protein  
interaction



**Figure 8.** A Map of TASS protein interactions. Nodes represent proteins while edges represent an interaction (25). Subunits of the TASS used as baits are highlighted in red. The inner, full circle represents interactions shared among the subunits. The broken outer circle represents interactions distinct to a given subunit. Transported effectors may more likely be found in the inner circle as they would interact with more than one subunit.

## Microarray



# Progress: Databases

- Updated PATRIC database schema to GUS 3.5 for all database instances
- Various improvements and additions to database infrastructure



# Curation: Two Parallel Approaches

## Bottom-Up Approach:

- High-Throughput Curation
- Systematic Analysis
- Standardized Pipeline

## Top-Down Approach:

- Targeted Analysis
- Specialized Curation
- Specialized Methods

**Diagnostics,  
Therapeutics &  
Vaccines**

**Target Discovery**

**Analysis**

**Curation**

50%

50%

# Top-Down Approach: Targeted Discovery of Countermeasure Candidates

- Community involvement
  - Goals defined by end users (organism experts and community they represent)
  - Experts' involvement at all stages
  - Commitment from end users for validations/follow-up
- Problem-driven curation and analysis

# Top down approach implementation: “Special Projects”

- *Brucella*
- *Rickettsia*
- Lyssavirus



# *Brucella* Special Project

- Project Goals
  - Identify virulence factors by comparative genomics
  - Identify functional polymorphisms
  - Obtain diagnostic markers
  - Incorporate methodology into genome annotation pipeline
- Collaborators
  - Stephen Boyle (PATRIC Organism Expert, VirginiaTech)
  - Yongqun “Oliver” He (University of Michigan)

some preliminary results

# Newly-Identified Genes in published *Brucella* genomes

“Newly-identified genes”:

- Genomic regions not previously associated with coding sequences
- ...which have full-length alignments with known or predicted proteins in closely related genomes

	Newly-identified Genes	With assigned function
<i>B. suis</i> 1330	50	14
<i>B. abortus</i> 2308	129	14

# Virulent vs. Attenuated Strain Comparison

- VBI has access to two unpublished *Brucella abortus* strains: 9-941 and S19
  - 9-941: parent strain
  - S19: attenuated mutant of 9-941 that developed spontaneously and is used as vaccine strain

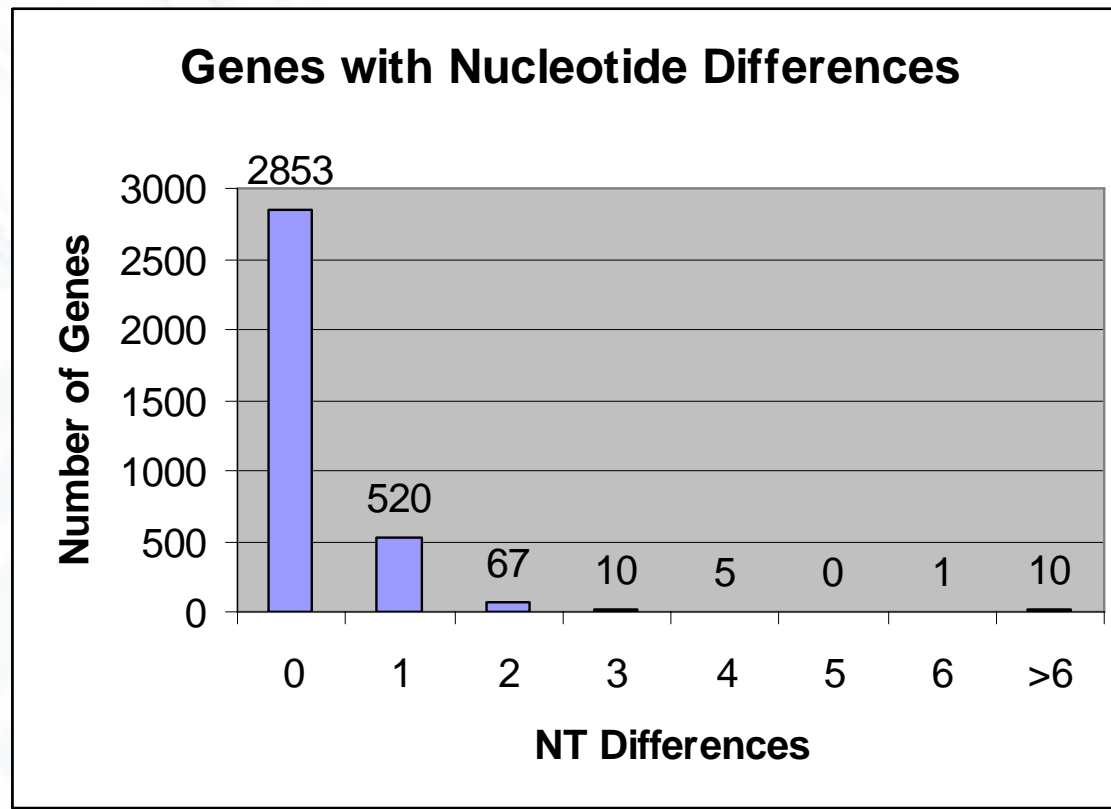
# Frameshifts and Premature Stops in *Brucella abortus* S19 ORFs

S19 Analysis	Size	Premature Stop	Frameshift	Both
Chr I	2.12 Mb	34	477	10
Chr II	1.16 Mb	31	320	7

In progress: filter to distinguish

- 1) sequencing errors
- 2) mutations shared with 9-941
- 3) mutations specific to S19

# Polymorphisms (Substitutions and Indels) between Orthologs in S19 and 9-941

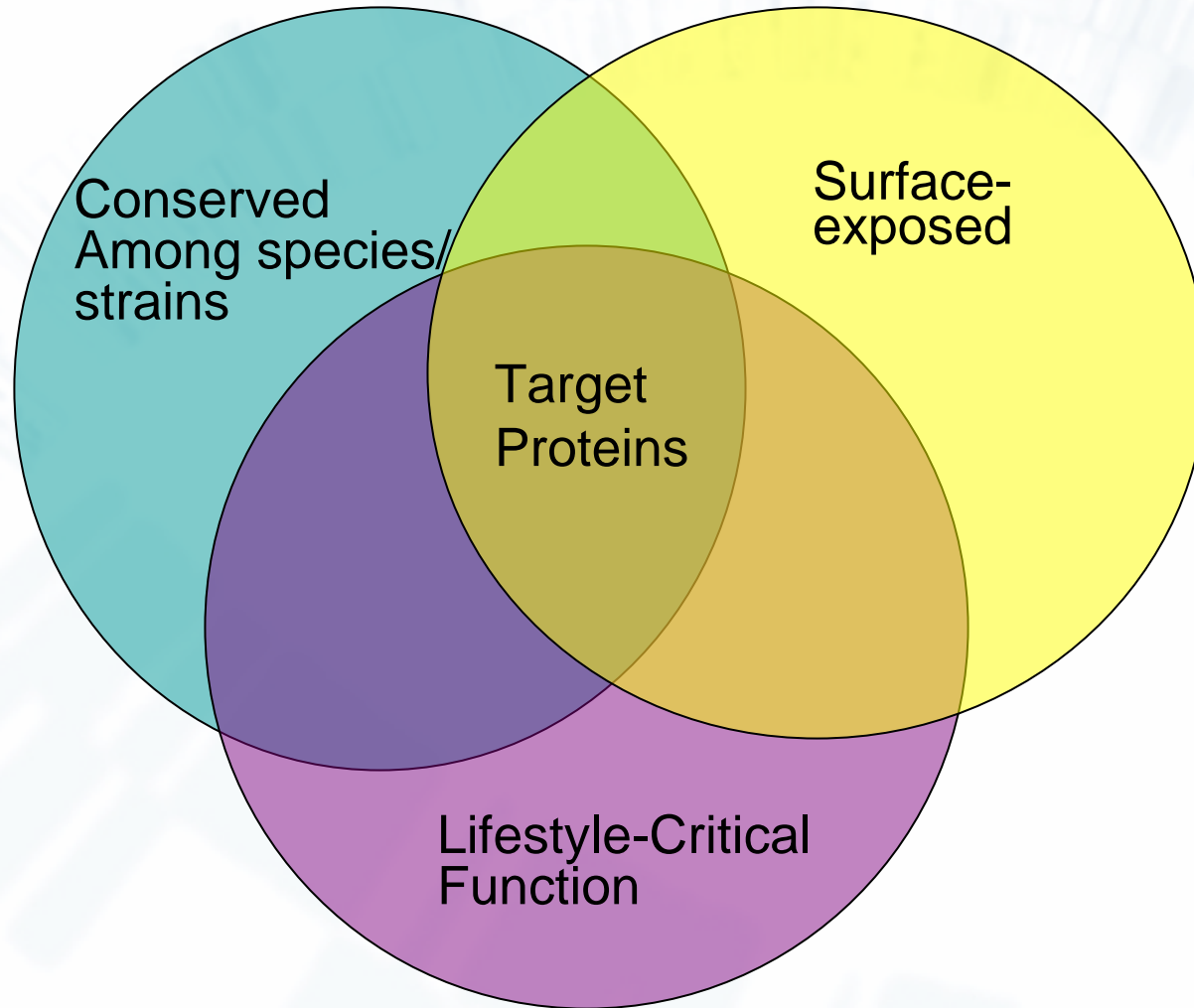


# ***Rickettsia* Special Project**

- Collaborator: Abdu Azad (PATRIC Org. Expert, U. of Maryland)



# Finding Good Vaccine Targets



# *Rickettsia prowazekii* Proteome

834

Signal-P  
Positive ANN?  
Positive HMM?

67

Secretome P

+

BOMP

Beta-Barrel?

9

Conserved?

6

Good Epitope?

Targets

IEDB

# Vaccine Target Candidate Proteins

- The six genes are annotated as **hypothetical proteins**
- Highly conserved among 7 annotated Rickettsial genomes
- Further analysis: the coded proteins seem to share common features/ motifs characteristic of **autotransporters, competence factors and secretion apparatus** proteins

# Lyssavirus Special Project

- Goal: develop an automated **phylogenetic classification system** (phylotyping) for Lyssaviruses
- Collaborators
  - Charles Rupprecht (PATRIC OE, CDC)
  - Alan Dickerman (VBI)

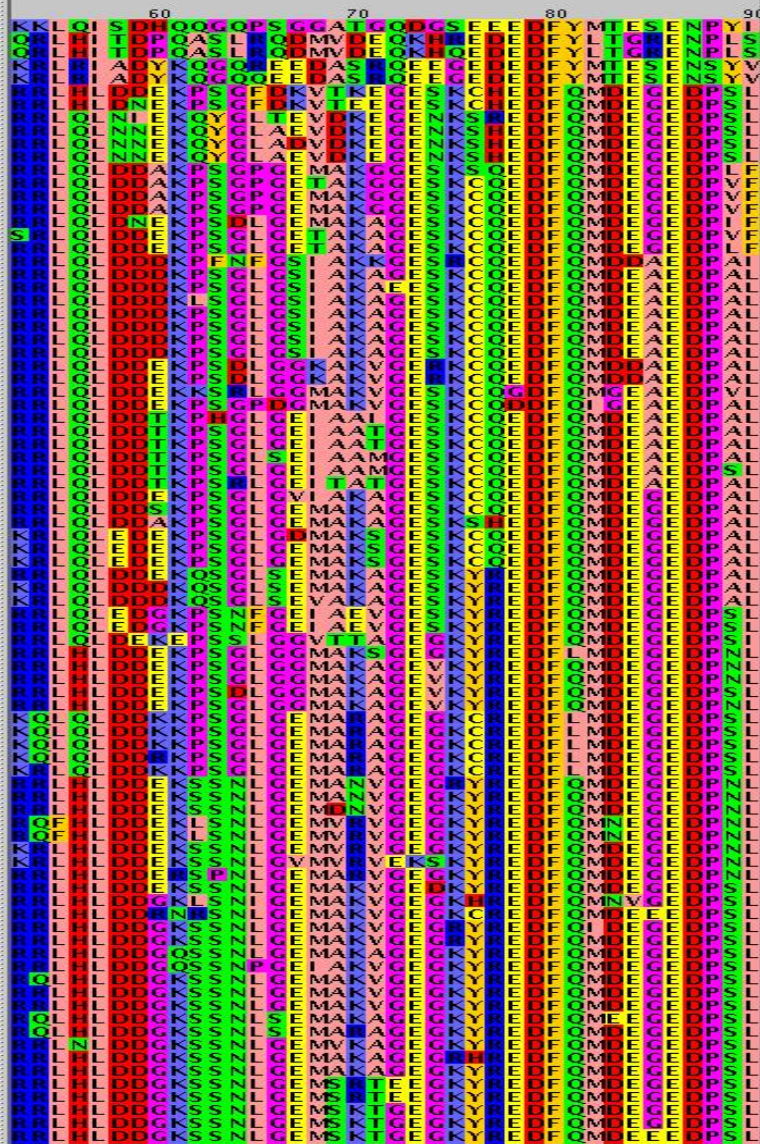
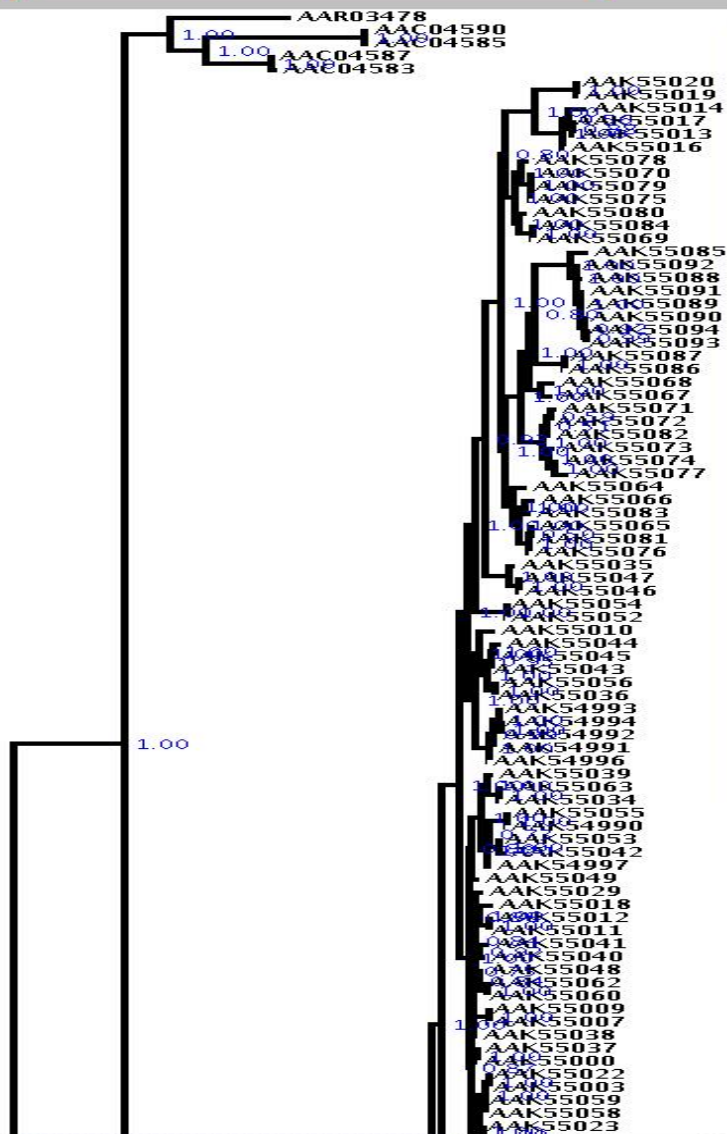
# Phylotyping

- Tool that will allow submission of sequences for multiple sequence alignment and generation of phylogenetic tree (web-based)
- Methods: MUSCLE + Mr.Bayes + VBI programs
- Tree Viewer juxtaposes tree view with MSA



File Tree Labels Rooting Debug

Tree Height = 2



# General plans for near future

- Steady state curation
  - Next release in June 2006
- Various improvements to infrastructure (including “automated curation”)
- LANL sequence annotation
- Special projects



# Organization of Relevant Meetings by VBI

- International Symposium on the Comparative Biology of Alpha-Proteobacteria
  - **April 26-29**, Blacksburg
  - (*Brucella* and *Rickettsia* are alphas)
- Computational Genomics '06
  - **October 28-31** (tentative), Baltimore
  - 3 main themes
    - **Infectious diseases**, automated annotation, and biological networks

# PATRIC

## PIs

B. Sobral  
J. Setubal

## Executive Committee

O. Crasta  
M. Czar (project management)  
R. Kenyon (project management)  
A. Purkayastha (curation)  
E. Snyder (bioinformatics)  
R. Will (software)

## Curators

C. Dharmanolla (literature)  
V. Dongre (Hep E)  
M. Hance (Rickettsia and Lyssavirus)  
D. Jukneliene (Coronavirus)  
L. Mackasmiel (Calicivirus)  
J. Shallom (Coxiella and Hep A)  
G. Yu (Brucella)

## Software Developers

N. Kampanya (web design and visualizations)  
J. Lu (database architect and administrator)  
M. Shukla (genome curation interfaces)  
J. Soneja (annotation pipelines)  
F. Zhang (lit. curation, web navigation and queries)

## Organism Experts

A. Azad (U. Maryland, Coxiella and Rickettsia)  
S. Baker (Loyola U., Coronavirus)  
S. Boyle (VT, Brucella)  
Y. He (U. Michigan, Brucella)  
Y. Khudyakov (CDC, Hep A)  
XJ Meng (VT, Hep E)  
C. Rupprecht (CDC, Lyssavirus)  
J. Vinje (CDC, Calicivirus)

## Collaborators

J. Gabbard and D. Hix (VT, usability engineering)  
N. Ramakrishnan (VT, data mining)

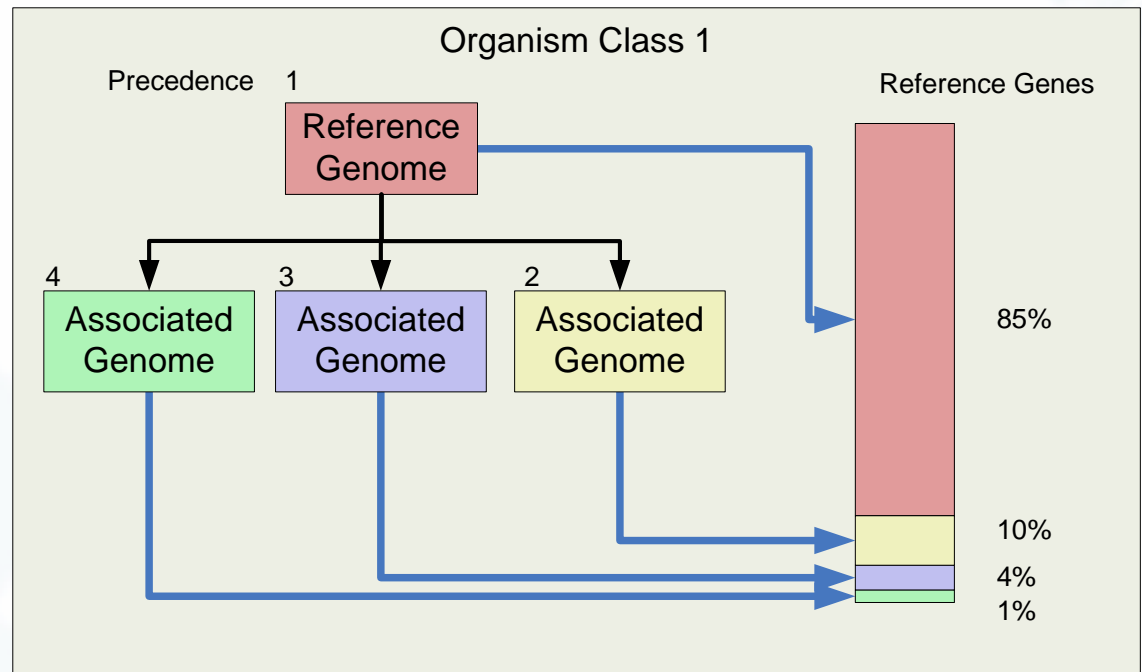


# Thank You

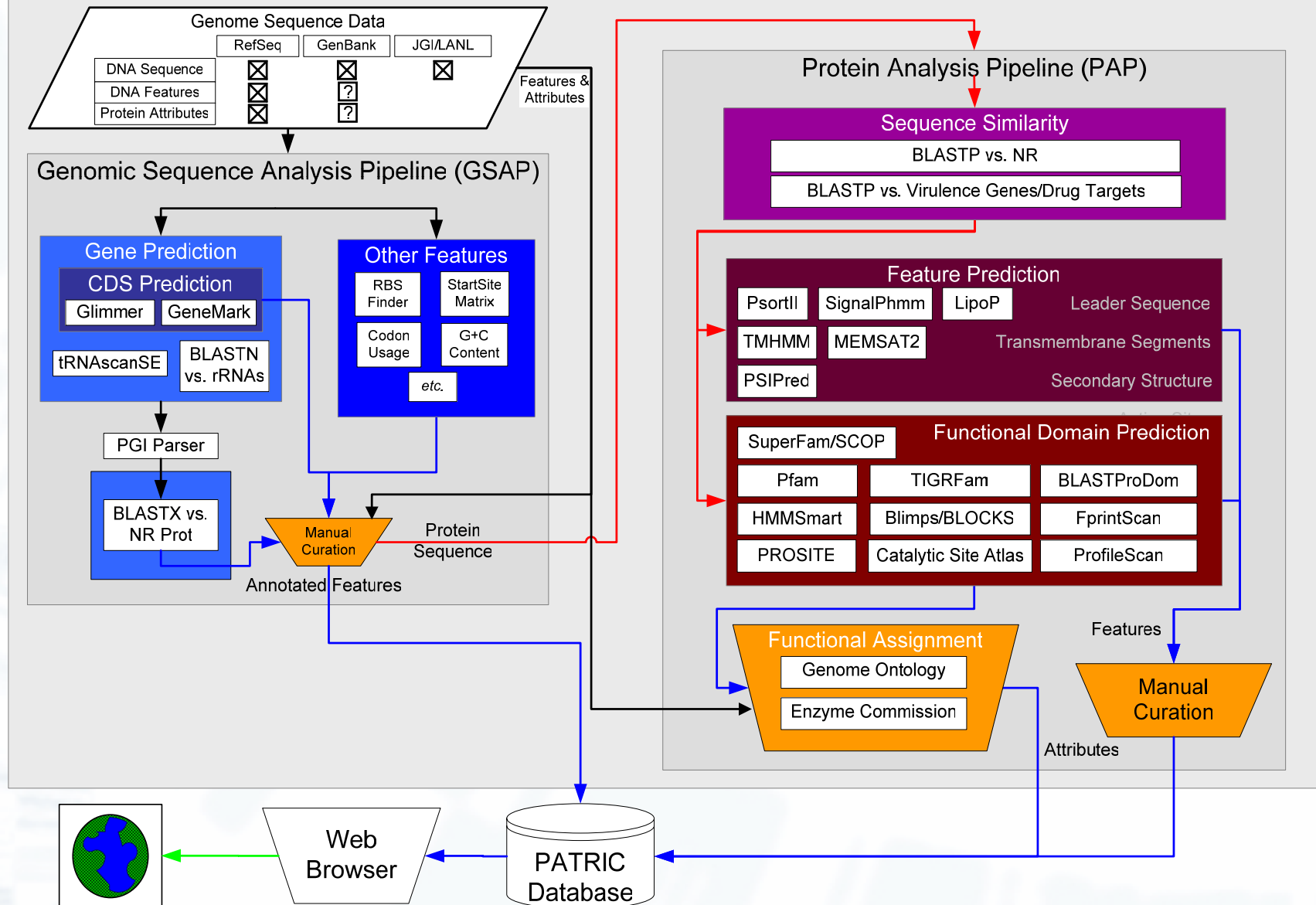
# What is a Reference Genome?

- Genome sequence of the type species of a genus.
- A lab strain that has been extensively characterized.
- A representative of a phylogenetic subgroup (community choice)

**The Reference Genome annotation is the starting point to create a Reference Gene Set for that Pathosystem.**



# Genome Analysis Pipeline (GAP)



# Reference Genomes Curated for the Dec. 22, 2005 Release.

Pathosystem	Complete Genomes	Genome Size	Reference Genomes	DNA curation	Protein Curation
<b>Bacteria</b>					
<i>Brucella</i>	4	3.2 Mb	1	Y	auto
<i>Coxiella</i>	1	1.9 Mb	1	Y	auto
<i>Rickettsia</i>	7	1.2 Mb	1	Y	auto
<b>Viruses</b>					
Calicivirus	64	7.5 kb	13	Y	N
Coronavirus	163	31 kb	16	Y	N
Hepatitis A	16	7.4 kb	1	Y	N
Hepatitis E	48	7.2 kb	5	Y	N
Lyssavirus	12	11.9 kb	1	Y	N

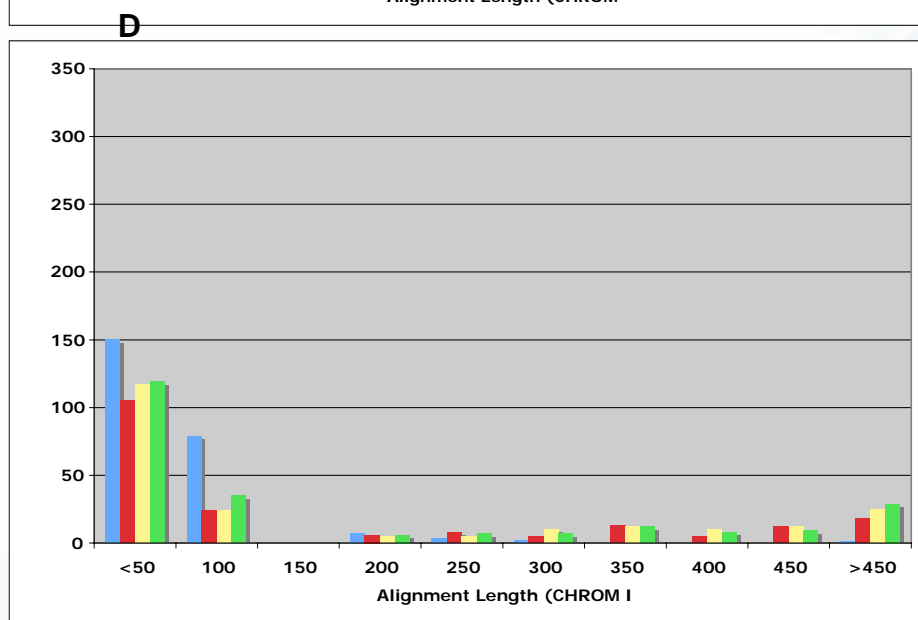
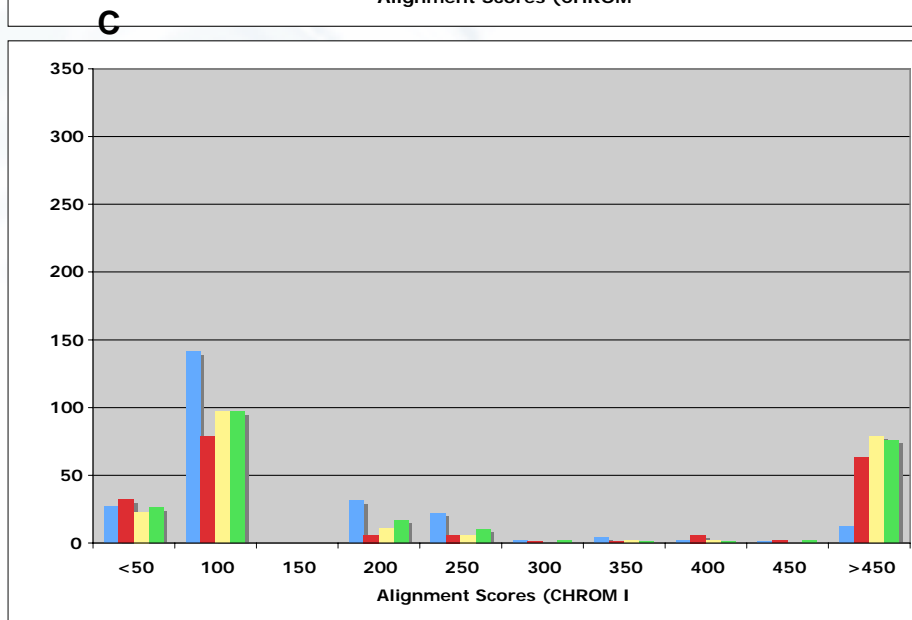
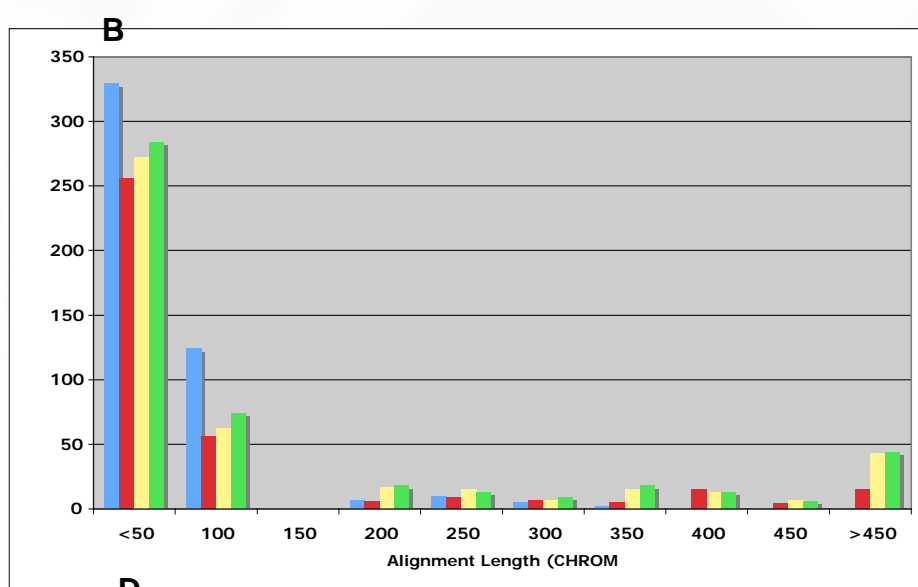
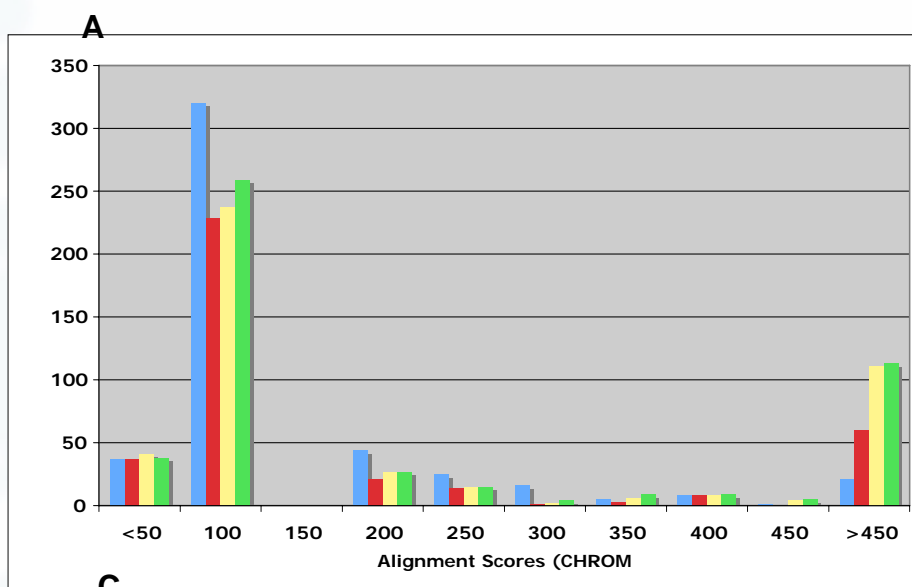
# Pathway Analysis: plan

- Metabolic pathway prediction/visualization hyperlinked to protein annotation pages (Pathway Tools, P. Karp)
- Signaling pathways, protein complexes, gene regulation pathways also to be incorporated
- Host/Pathogen interaction visualizations in development
- All pathway proteins hyperlinked to proteomic experimental data

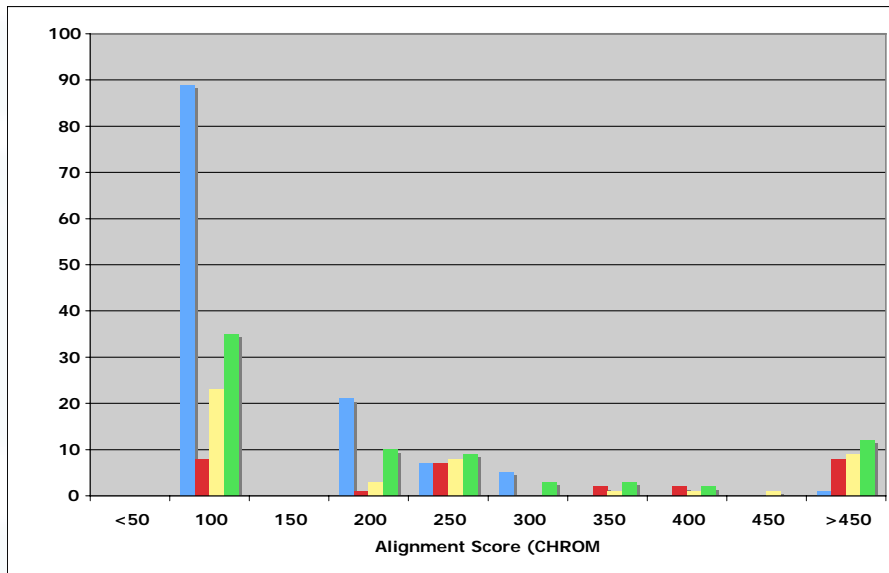
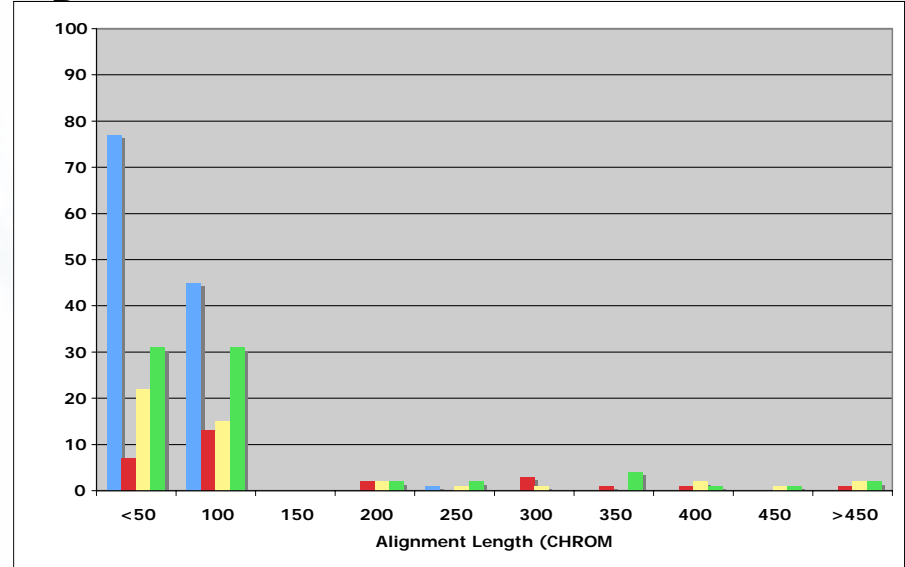
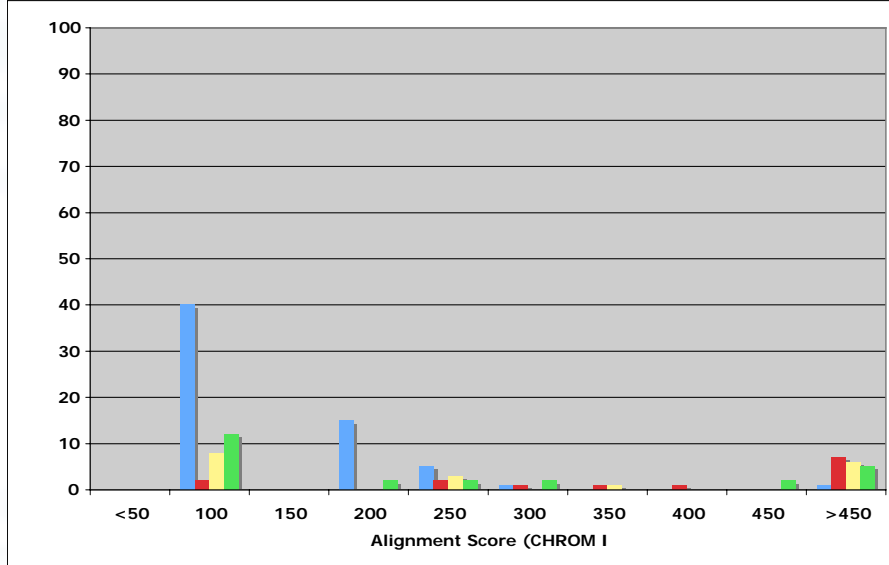
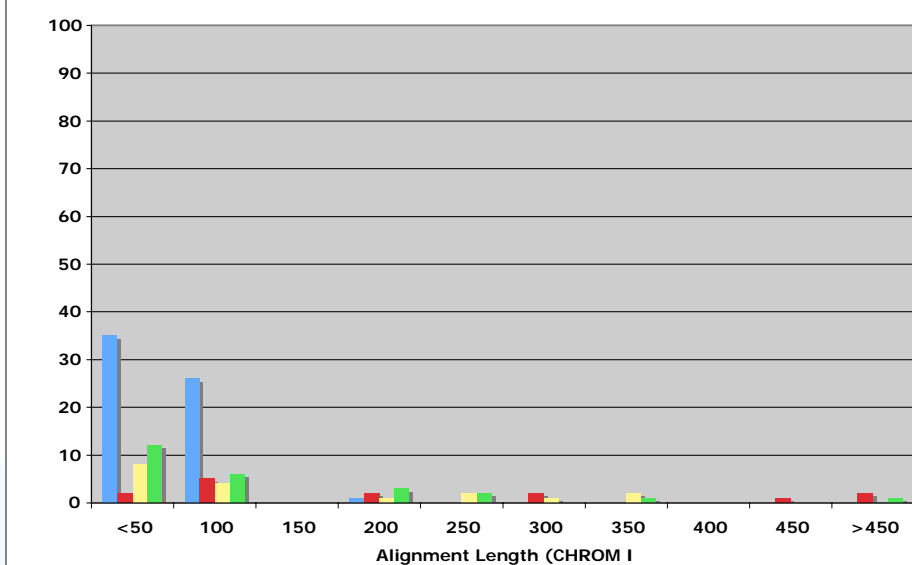


# A Multi-Task Procedure for Critical Analysis of Sequence Variants and Comparative Genome Analysis

1. To identify missed genes, genes that are involved in frame-shifts, in-frame stop codons, insertions, and deletions
2. To comparatively analyze these gene variants among closely-related genomes to identify genome-specific genetic factors
3. Build as a part of our automatic genome annotation pipeline and offer as open source software for research community.

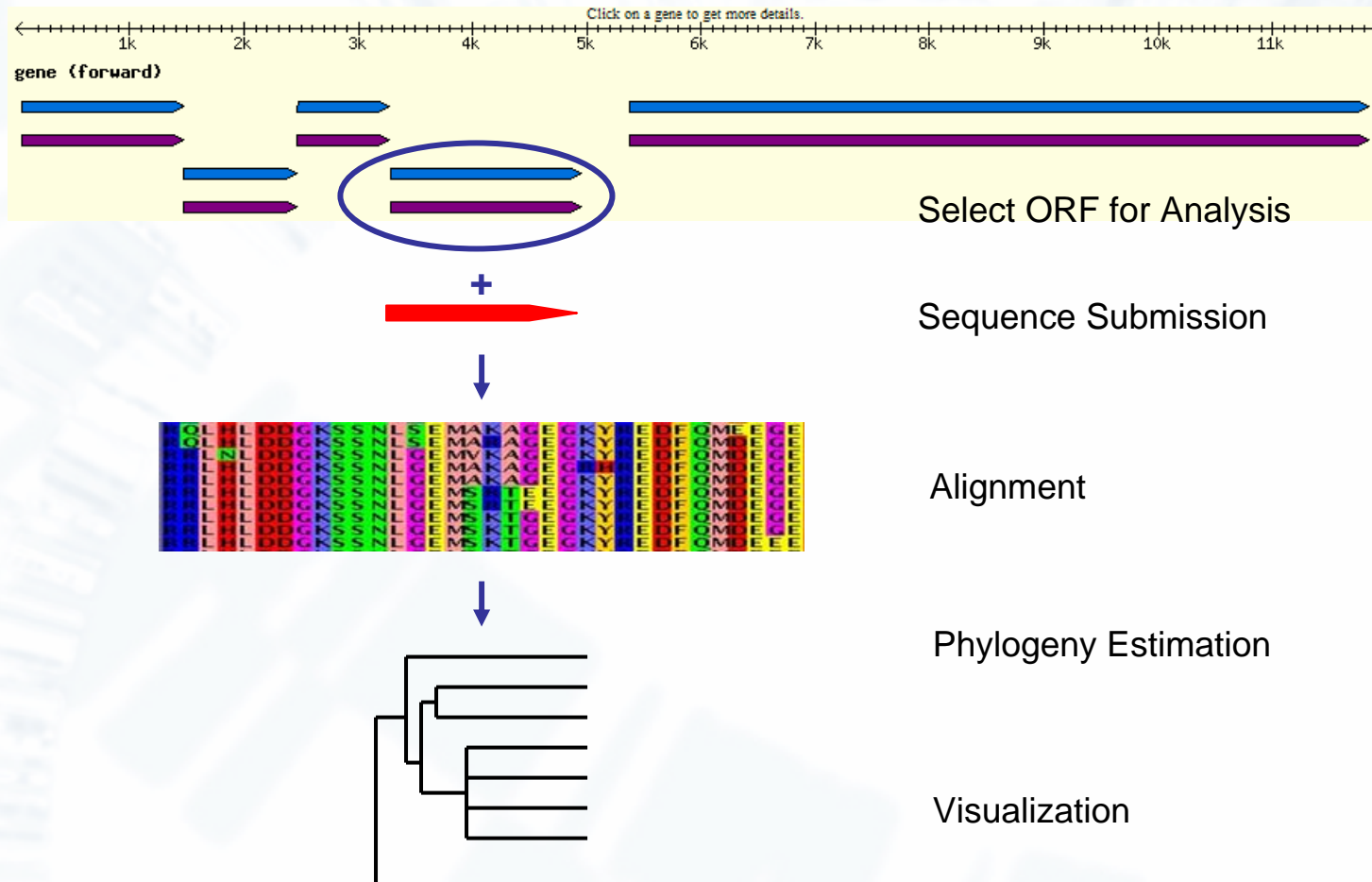


Distribution of homologs revealed in the inter-genic regions of four *Brucella* genomes: *Brucella melitensis* 16M (Blue), *Brucella suis* 1330 (red), *Brucella abortus* str. 9-941 (yellow) and *Brucella abortus* str. 2308

**A****B****C****D**

Distribution of missed genes revealed in the inter-genic regions of four *Brucella* genomes:  
*Brucella melitensis* 16M (Blue), *Brucella suis* 1330 (red), *Brucella abortus* str. 9-941 (yellow)  
 and *Brucella abortus* str. 2308 (green)

# Phylotyping Application Workflow



# General plans for near future

- Steady state curation
  - Next release in June 2006
    - Nucleotide-level automated curation for all genomes
    - Automated protein-level for all reference genomes
    - Manual protein-level curation of one bacterial genome
- Various improvements to infrastructure
- LANL sequence annotation
- Special projects

# PATRIC Representation in Community Meetings

Meeting	Month	Location
Rickettsia	June 2005	Spain
Nidovirus (incl. Coronaviruses)	June 2005	Colorado
Int'l Union of Microbiological Societies	July 2005	California
Int'l Virus Database Meeting	June 2005	Missouri
Brucellosis	October 2005	Mexico
Rabies	October 2005	Canada
Genome Informatics	October 2005	CSHL, NY
Biocurator	December 2005	California



# Meeting Activities

- Abstracts
  - 4 Posters
  - 4 Presentations
- Sponsored Booth at IUMS
  - Cosponsored with VBRC
- Web-based questionnaires for four meetings
  - To gauge the bioinformatic needs of community
  - Nidovirus: 18 responses
  - IUMS: 9 responses
  - Rabies: 3 responses
  - Brucellosis: 0 responses

# Survey Responses

Resource	IUMS Response No.	IUMS Mean Score	Nidovirus Response No.	Nidovirus Mean Score
Universal Primers	7	2.7	13	1.8
Epitope Database	7	2.7	14	2.4
Viral gene expression and assembly pathways	8	2.8	14	2.4
Phylogenetic trees	7	3.0	13	2.0
Comparative Genomics	7	3.0	12	2.3
Affect of viral proteins on host pathways	8	3.1	14	1.6
Modelling of active sites of enzymes	7	3.3	12	2.4
High-throughput data	7	3.3	12	2.5
Epidemiological data	7	3.3	13	2.5
Literature collection	7	3.3	14	1.6
3D structural data	8	3.5	14	1.9

1=Highest Priority

5=Lowest Priority

# Follow Up to Survey

- Difficult to gauge needs
  - Low response to questionnaires
- Different needs of different communities
- Gap between bioinformatics resources and workflow of countermeasures development
- Approach:
  - Expedite response to one community through Special Projects; results directly applicable to one community
  - Resultant Use-Case should help build broadly applicable analysis resources

# SWG meetings

- Face-to-face June 2005
- Conference call January 2006
- Next one: face-to-face June 2006